

ANÁLISIS COMPARATIVO DE TÉCNICAS ESTADÍSTICAS Y DE APRENDIZAJE PARA EVALUAR LA SUSCEPTIBILIDAD DEL TERRENO A LOS DESLIZAMIENTOS SUPERFICIALES EN EL PIRINEO CATALÁN.

**Samuel AMORIM^{1*}, Jordi COROMINAS¹, Nieves LANTADA¹, Cristina BAEZA¹,
Modesto PORTILLA¹ y Cecilio ANGULO²**

¹ Departamento de Ingeniería del Terreno, Cartográfica y Geofísica (ETCG)
Universidad Politécnica de Cataluña (UPC)

² Ingeniería de Sistemas, Automática e Informática Industrial (ESAI)
Universidad Politécnica de Cataluña (UPC)

RESUMEN

Este documento presenta la comparación de tres metodologías de tratamiento de datos para la evaluación de la susceptibilidad a la ocurrencia de deslizamientos superficiales. Por un lado dos modelos estadísticos como el análisis discriminante (AD) y la regresión logística (RL) y por otro el modelo de aprendizaje de las redes neuronales artificiales (RN). Se ha obtenido la probabilidad espacial de rotura para los diferentes métodos partiendo de una base de datos idéntica referente a la misma área. En la validación y comparación de los modelos se utilizaron los histogramas de frecuencia de los grupos estables e inestables, índices de densidad relativa de deslizamiento, curvas de tipo ROC, tasa de éxito, predicción y porcentaje acumulado. Se concluye que los modelos son semejantes, no obstante, el AD presentó una capacidad predictiva superior en el índice de densidad relativa y la RN resultados ligeramente mejores en las demás evaluaciones.

1. INTRODUCCIÓN

La evaluación de la susceptibilidad es el primer paso del análisis de riesgo. La determinación de la susceptibilidad en términos de una probabilidad espacial permite, además de la comparación directa entre los distintos métodos de predicción, la evaluación cuantitativa de la peligrosidad y del riesgo.

La predicción de áreas potenciales de rotura presenta una gran incertidumbre debido a los pocos datos disponibles y las deficiencias de los modelos de susceptibilidad (Zêzere, 2002). Es necesario el desarrollo de metodologías rápidas y fiables que permitan una primera aproximación y delimitación de las áreas más susceptibles. Estas metodologías deben ser, en la medida de lo posible, suficientemente generalizables de forma que permitan su aplicación a otras zonas semejantes con una mínima adaptación.

Las metodologías de evaluación de susceptibilidad suelen dividirse en cualitativas y cuantitativas. Las clasificaciones propuestas son subjetivas y dependen de los aspectos considerados por cada autor. Soeters y van Westen (1996) distinguen: los métodos heurísticos que se dividen en análisis geomorfológico (determinación directa en campo de la peligrosidad) y la combinación cualitativa de mapas (a partir de opinión experta); los métodos estadísticos bivariante y multivariante y finalmente los métodos determinísticos que evalúan el factor de seguridad mediante el análisis de estabilidad de la laderas teniendo en cuenta los parámetros de resistencia del terreno. Recientemente se están aplicando en el análisis de susceptibilidad diferentes técnicas basadas en el aprendizaje y en las funciones de pertenencia (redes neuronales artificiales y lógica difusa), dando origen también a métodos híbridos con la combinación de algunos de los métodos citados (Ayalew et al., 2005).

Los métodos estadísticos según Carrara et al. (2008), se basan en el análisis de las relaciones entre la distribución espacial de factores condicionantes de la inestabilidad de la ladera y la de los deslizamientos observados. Se asume que los factores que causaron roturas en una región específica son similares a los que podrían generar deslizamientos en el futuro. Utilizando una misma hipótesis, las técnicas de aprendizaje buscan definir un modelo que clasifique determinados patrones o ajuste una función de interés minimizando numéricamente el error encontrado entre la salida estimada y la respuesta disponible de un conjunto de datos previamente conocido. En esta comunicación se presenta la aplicación y comparación de dos modelos estadísticos como el análisis discriminante (AD) y la regresión logística (RL) y el modelo de aprendizaje de redes neuronales artificiales (RN).

2. ÁREA DE ESTUDIO

La región seleccionada como área de estudio tiene una extensión aproximada de 40km², con dimensiones de 8x5km, está situada en el PrePirineo Oriental, más específicamente en las proximidades del municipio de la Pobla de Lillet, en la comarca catalana del Berguedà.

El substrato rocoso de la región es de carácter eminentemente carbonatado con edades que van desde el Devónico medio hasta el Eoceno medio. Se distribuye en diferentes unidades estructurales (Unidad del Cadí y manto del Pedraforca) caracterizadas por una tectónica alpina, y ocupan franjas del relieve dispuestas en sentido E-W (Santacana et al., 2002).

Los deslizamientos superficiales se consideran aquellos cuya superficie de rotura se sitúa

a una profundidad media de 1 a 2 metros, afectando en general la formación superficial que recubre la ladera y/o la parte alterada del sustrato cuando este es de naturaleza lutítica (Santacana, 2001). En la región de estudio, los deslizamientos superficiales afectan mayoritariamente a formaciones coluviales cuaternarias y, en menor medida, a niveles meteorizados de las arcillitas del Permotrías, arcillitas del Cretácico superior y turbiditas del Eoceno. (Santacana et al., 2002).

La base de datos que utilizamos es la de Santacana (2001) a escala 1:5000. Contiene 280 deslizamientos producidos durante los aguaceros de noviembre de 1982. Contiene además parámetros obtenidos mediante fotointerpretación y reconocimiento de campo (espesor de la formación superficial, uso del suelo y cobertura vegetal) y variables referentes a la geometría de la ladera y de la cuenca obtenidas a partir del Modelo Digital de Elevaciones (MDE) de estructura raster con píxeles de 15x15m (Figura1).

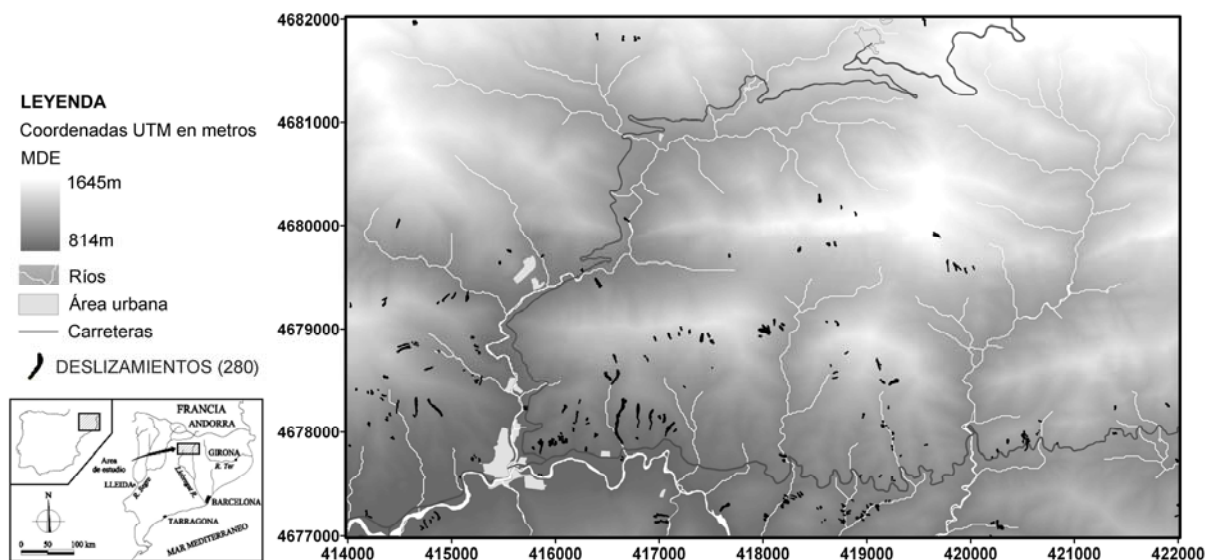


Figura 1. Mapa de situación y Modelo Digital de Elevaciones del área de estudio con el inventario de deslizamientos.

3. SELECCIÓN DE VARIABLES

Gran parte de los métodos de evaluación de la susceptibilidad se basan en el análisis de los parámetros geológicos y geomorfológicos del terreno y sus contribuciones sobre las condiciones de rotura y sobre la movilidad de los deslizamientos superficiales. Las variables utilizadas para el análisis de susceptibilidad, han sido seleccionadas teniendo en cuenta su relación con la aparición de roturas en la zona de estudio. Santacana (2001) partió de un conjunto de 16 variables condicionantes reducidas posteriormente a 12 mediante análisis factorial de componentes principales, eliminando las variables dependientes de menor

significancia estadística. Posteriormente, la sucesiva aplicación de diferentes modelos de análisis discriminante con inclusión de variables por pasos definió un subconjunto de 7 variables que explicaban de mejor forma la susceptibilidad. Las variables obtenidas de esta manera, utilizadas en este estudio son: uso del suelo y cubierta vegetal (VEG), espesor de la formación superficial (GROSOR), altitud sobre el nivel del mar (DTMFILL), pendiente de la ladera (PEND), grado de concavidad/convexidad del relieve en sección transversal a la pendiente (PLA), longitud máxima ponderada de cuenca vertiente (LLONG) y pendiente media de la cuenca (PENDM). Las relaciones entre cada factor y los deslizamientos así como los detalles de su definición, cálculo y proceso de selección se pueden encontrar en Santacana (2001) y Santacana et al. (2002).

4. EVALUACIÓN DE LA SUSCEPTIBILIDAD

Los dos modelos estadísticos han sido implementados con ayuda del programa SPSS mientras que el de aprendizaje con el programa MATLAB. Los datos de entrada y los resultados fueron elaborados, gestionados y analizados utilizando el programa ArcGIS.

El pretratamiento de datos fue común para todos los modelos y constó de división aleatoria del conjunto de celdas inestables (centroide de la zona de ruptura de los deslizamientos inventariados) en muestras iguales para entrenamiento y validación, generación de celdas estables (elegidos aleatoriamente entre todos los demás píxeles del mapa raster), cambio de escala de covariables y mezcla de casos inestables y estables. La verificación de ajuste de las variables a una distribución normal y su consiguiente transformación, si bien es innecesaria para el análisis con Regresión Logística y Redes Neuronales, se realizó para todos los modelos.

4.1. Análisis Discriminante

El análisis discriminante es una técnica estadística para clasificar individuos u objetos en grupos mutuamente excluyentes basándose en un conjunto de variables independientes (Carrara et al., 2008). En su aplicación a la rotura de laderas, las variables independientes están constituidas por los factores condicionantes de la inestabilidad. La variable categórica dependiente es un factor de agrupación que coloca cada individuo de la muestra en uno y sólo uno de los grupos definidos a priori, en base a una relación lineal (Dillon y Goldstein, 1984). El análisis discriminante permite la clasificación en dos o más grupos, aunque en nuestro caso fue aplicado para distinguir entre estable e inestable.

Una vez calculadas la función discriminante, a cada objeto se le puede asignar una puntuación (*score*) o valor de la función discriminante dado por la variable dependiente Y (Ecuación 1), donde C son los coeficientes de clasificación estimados y X son las variables independientes más importantes estadísticamente para el proceso discriminante.

$$Y = C_0 + C_1X_1 + \dots + C_nX_n \quad (1)$$

El programa SPSS aplica la regla de Bayes para la clasificación de los casos a partir de los valores de la función discriminante, obteniendo la probabilidad de que el objeto pertenezca a uno u otro grupo establecido (Visauta y Martori, 2003).

4.2. Regresión Logística

La regresión logística es útil para los casos en los que se desea predecir la presencia o ausencia de una característica o resultado según los valores de un conjunto de indicadores (Lee et al., 2007). La relación multivariante de regresión que permite establecer la relación espacial entre ocurrencia de deslizamientos y sus factores condicionantes se expresa por la Ecuación 2, también lineal:

$$Z = \beta_0 + \beta_1X_1 + \dots + \beta_nX_n \quad (2)$$

Donde Z es la variable dependiente que expresa la presencia (1) o ausencia (0) de deslizamientos, β_0 es el intercepto, $\beta_1 \dots \beta_n$ son los coeficientes estimados que miden la contribución de las variables independientes ($X_1 \dots X_n$) definiendo los factores que más afectan la estabilidad (Duman et al., 2006).

La principal ventaja de la regresión logística es que por medio de la adición de una función de enlace, permite que las variables sean continuas, discretas o una combinación de ambas sin que necesariamente deban tener distribuciones normales (Lee et al., 2007). Esta característica amplía su aplicabilidad a un rango más amplio de situaciones que el análisis discriminante. La función logística de enlace es aplicable a los casos en que la variable dependiente es dicotómica (Atkinson y Massari, 1998). De ella deriva la probabilidad (p) de que la variable dependiente (Z) tenga valor 1 definida por la Ecuación 3.

$$p_i = \frac{1}{1 + e^{-Z_i}} \quad (3)$$

4.3. Redes Neuronales

Las redes neuronales artificiales son técnicas de la inteligencia artificial basadas en el comportamiento del cerebro humano, el aprendizaje y la posterior operación de los procesos. Su propósito es construir un modelo de generación de datos del proceso de tal forma que la red pueda generalizar y predecir salidas para datos de entrada de los cuales no tiene información (Lee et al. 2007).

Se puede definir la RN utilizada como un perceptrón multicapa (MLP, *MultiLayer Perceptron*) de propagación hacia delante (*feedforward*) con siete entradas (factores

condicionantes) cuatro neuronas en la capa oculta y una salida (susceptibilidad de ocurrencia de deslizamientos). El algoritmo de entrenamiento supervisado utilizado fue el de retropropagación del error (*backpropagation*), modificado por Levenberg-Marquardt (Demuth et al., 2006).

El entrenamiento supervisado consiste en la optimización de la estructura y/o parámetros del modelo con el propósito de minimizar una función de error con base en los resultados obtenidos para el conjunto de datos cuyo resultado previo es conocido (Nelles, 2000).

En la Figura 2 se presenta la estructura de red utilizada en este trabajo con la nomenclatura vectorial del programa MATLAB. La capa 1 es la denominada capa oculta y la capa 2 la capa de salida. Las funciones de activación elegidas fueron la tangente hiperbólica en la capa oculta y la función logística (sigmoidea) en la capa de salida.

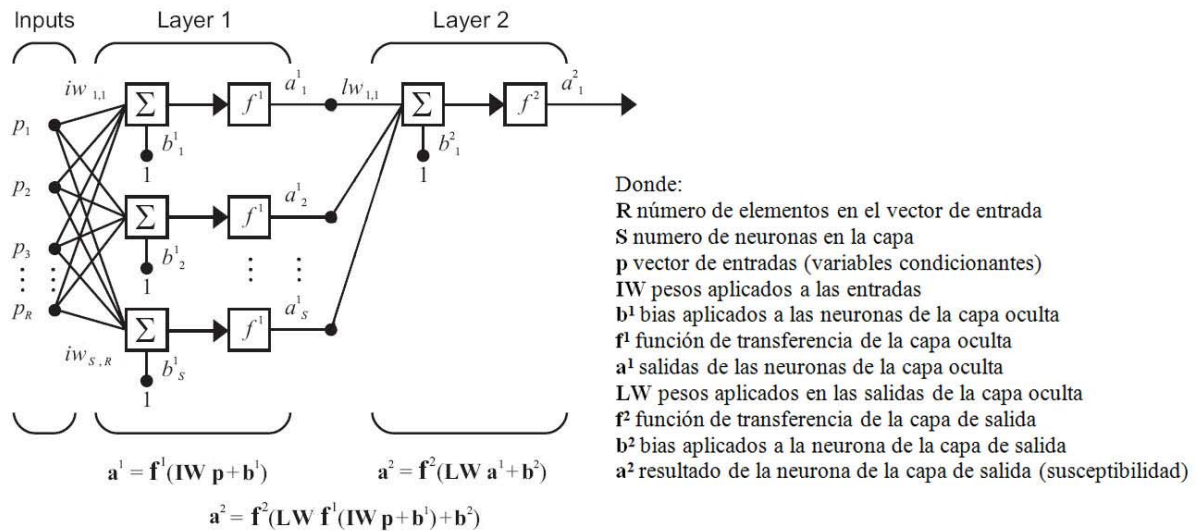


Figura 2. Perceptrón multicapa utilizado, esquema adaptado de Demuth et al. (2006)

Un perceptrón multicapa entrenado para clasificación utilizando *backpropagation* se aproxima a las funciones clasificadoras bayesianas más tradicionales. La salida de estos perceptrones multicapa se asemeja a la función de probabilidad posterior de las clases entrenadas (Ruck et al., 1990).

5. MAPAS DE SUSCEPTIBILIDAD Y ANÁLISIS COMPARATIVO

En este trabajo la susceptibilidad obtenida se presenta en cinco clases (muy baja, baja, media, alta y muy alta) definidas por intervalos iguales de la probabilidad de ocurrencia de deslizamientos. Este tipo de clasificación es indicada cuando las probabilidades encontradas por los distintos métodos son compatibles, permitiendo una comparación adecuada y objetiva

(Chung y Fabbri, 2003). Detalles de los mapas obtenidos con las diferentes metodologías pueden verse en la Figura 3. La Tabla 1 trae los coeficientes de los modelos de AD y RL. No se presentarán en este trabajo los pesos y bias obtenidos para las diferentes neuronas en las distintas capas del modelo RN por no ser directamente comparables con los pesos de AD y RL.

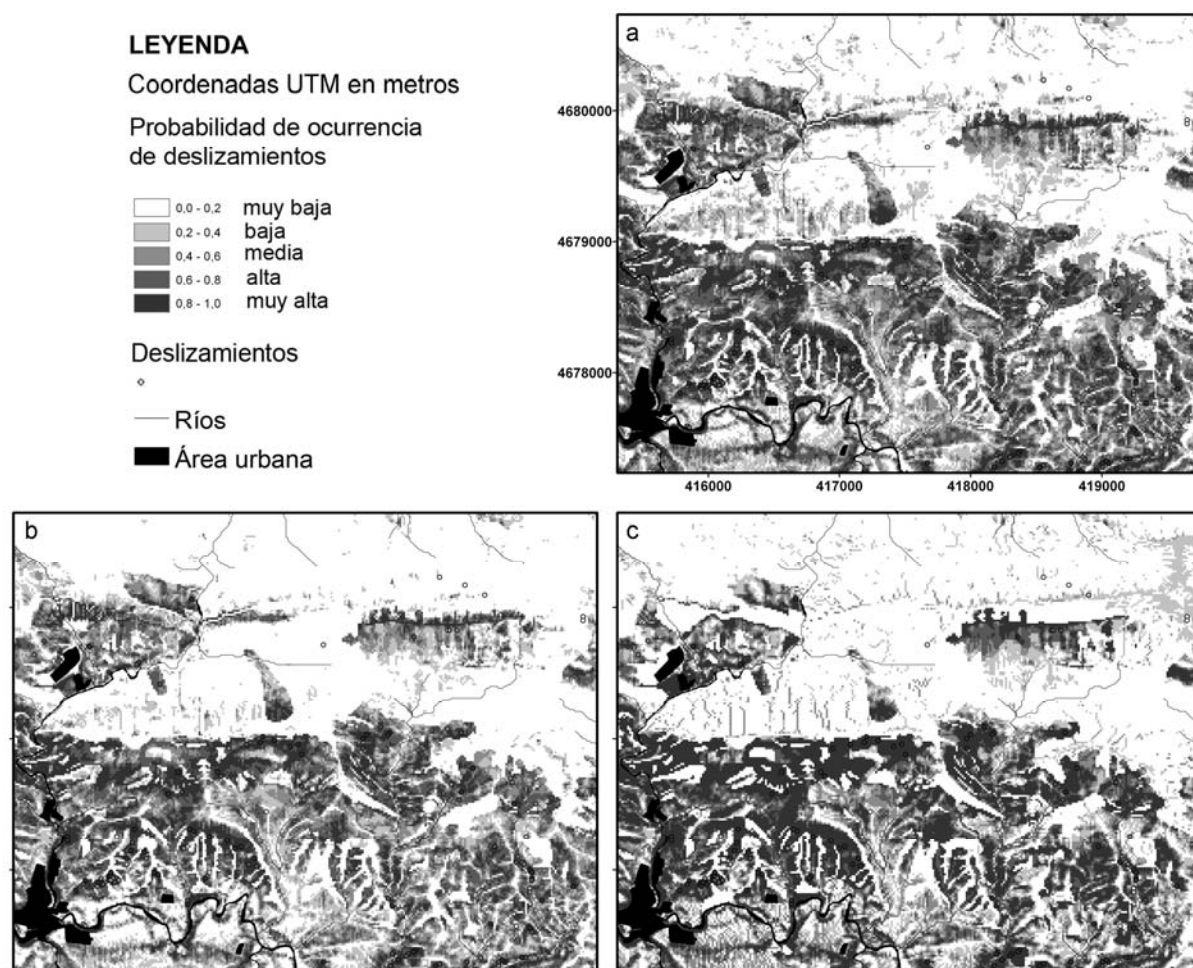


Figura 3. Detalle de los mapas de probabilidad de ocurrencia para modelos generados por a) Análisis Discriminante, b) Regresión Logística y c) Red Neuronal.

		Variables							Constantes
		DTMFILL	PLA	LLONG	PENDM	VEG	GROSOR	PEND	
AD	C_n	-1,666	-9,119	1,435	-2,360	-1,169	2,283	7,002	1,376
RL	β_n	-2,793	-14,689	2,696	-4,055	-2,485	4,524	12,014	1,090

Tabla 1. Coeficientes no tipificados de la función canónica discriminante (Ecuación 1) y de la relación multivariante lineal de la regresión logística (Ecuación 2) para las variables utilizadas (ver texto en Apartado 3).

Los histogramas de frecuencia de las poblaciones estables e inestables presentan diferencias en sus distribuciones y áreas de solape (Figura 4). Sin embargo, las distancias entre los centroides de estas poblaciones son similares. La RN presenta un área de solape menor y una distancia algo mayor entre las poblaciones (Tabla 2). A su vez en la RL destaca la gran frecuencia de celdas estables en la más baja susceptibilidad (Figura 4).

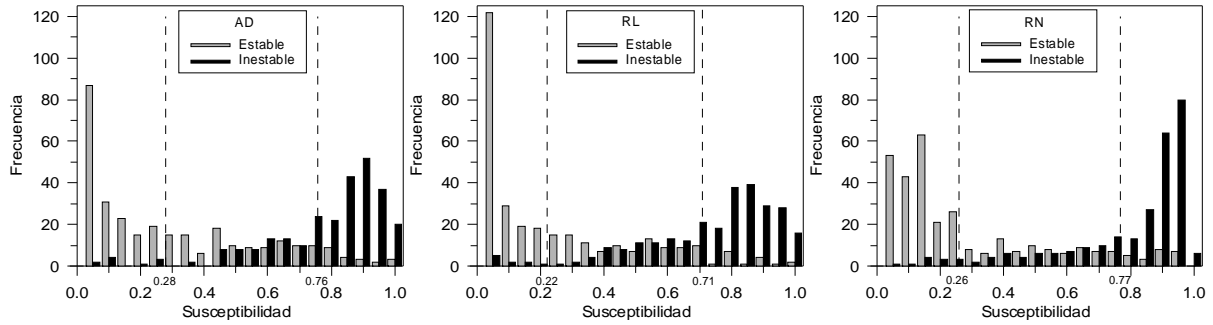


Figura 4. Histograma de frecuencia de grupos estables e inestables para diferentes modelos

Centroides	AD	RL	RN
estables	0,279	0,221	0,259
inestables	0,757	0,709	0,768
distancia	0,478	0,488	0,509

Tabla 2. Valor de los centroides de los grupos estables e inestables y distancia encontrada entre ellos para los diferentes modelos.

Para evaluar los resultados utilizamos el Índice de Densidad Relativa (IDR). El índice IDR se define en la Ecuación 4 por la proporción de roturas por clase de susceptibilidad, normalizado por la densidad global de roturas (Baeza y Corominas, 2001).

$$IDR_i = \left[(n_i / N_i) / \sum (n_i / N_i) \right] \cdot 100 \quad (4)$$

Donde n_i es el número de celdas con deslizamientos observados para la clase de susceptibilidad i y N_i es el número total de celdas de esta clase.

El índice se ha obtenido considerando 5 clases con intervalos iguales, si modificamos el número de clases o el rango de los intervalos el IDR también cambiará. Es deseable que los valores de IDR sean nulos para las clases de baja susceptibilidad y creciente desde las medias a las más altas susceptibilidades. La Figura 5 presenta los resultados encontrados por los distintos modelos, donde se puede observar un mejor comportamiento del modelo de AD.

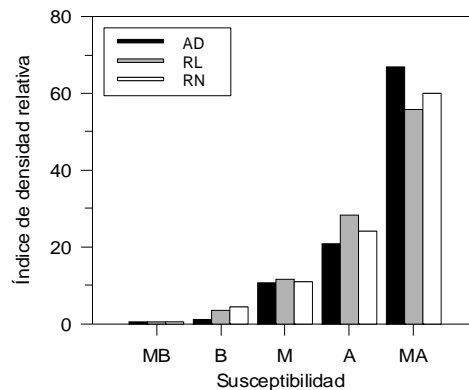


Figura 5. Índice de Densidad Relativa de Deslizamientos para los modelos evaluados.

Las curvas de porcentaje acumulado presentan separadamente los porcentajes de celdas con roturas observadas (celdas inestables) por intervalo discreto de la susceptibilidad y la extensión de área incluida en cada intervalo (celdas totales). Al definir las, Duman et al. (2006) analiza los valores encontrados para el punto de corte que clasifica las celdas dicotómicamente en estables o inestables, en este caso 0,5 (Figura 6). En el análisis de los resultados obtenidos para este mismo valor de corte se puede observar que en el AD casi el 90% de los deslizamientos ocurren en valores de susceptibilidad considerado inestables ($>0,5$), si bien el área ocupada por el conjunto de celdas consideradas inestables es del 28%. De todos modos, los valores obtenidos en los tres métodos son muy parecidos (Tabla 3).

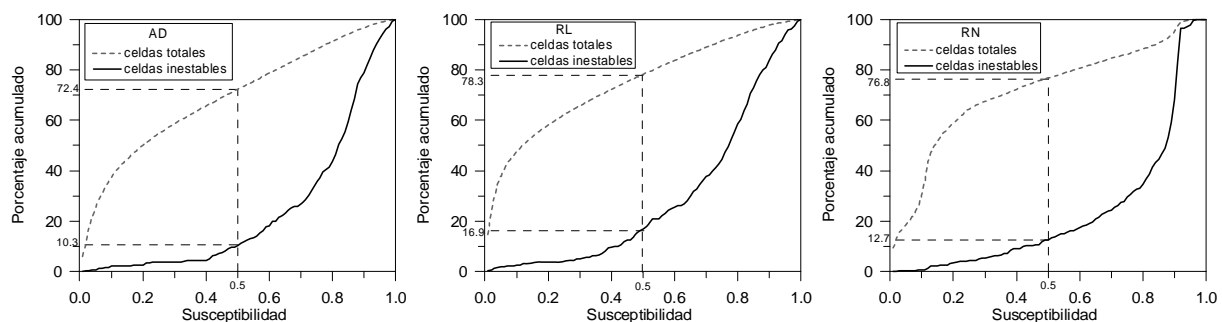


Figura 6. Curvas de Porcentaje Acumulado de celdas inestables observadas y celdas totales para cada modelo.

Porcentaje acumulado	AD	RL	RN
celdas inestables bien clasificadas	89,7	83,1	87,3
celdas totales clasificadas como inestables	27,6	21,7	23,2

Tabla 3. Porcentaje acumulado de celdas inestables y celdas totales clasificadas como inestables (valor de corte de 0,5) para cada modelo.

Adicionalmente, las curvas de porcentaje acumulado permiten analizar cualquier valor o

intervalo de susceptibilidad, lo que representa una ventaja frente al Índice de Densidad Relativa (IDR) calculado con rangos definidos.

Al establecer un punto de corte para evaluar la capacidad del modelo de discriminar entre dos grupos se compara los valores observados con los predichos, estableciendo los porcentajes de verdaderos positivos (VP), falsos positivos (FP), verdaderos negativos (VN) y falsos negativos (FN). La curva ROC (del acrónimo inglés *Receiver Operating Characteristic*) es la visualización grafica del análisis de todos los posibles puntos de corte y representa los pares (1-especificidad, sensibilidad) de cada uno de ellos (Carrara et al., 2008). Las Ecuaciones 5 y 6 definen la sensibilidad o fracción de verdaderos positivos (FVP) y la especificidad o fracción de verdaderos negativos (FVN).

$$\text{sensibilidad} = \text{FVP} = \text{VP}/P = \text{VP}/(\text{VP} + \text{FN}) \quad (5)$$

$$\text{especificidad} = \text{FVN} = 1 - \text{FFP} = \text{VN}/N = \text{VN}/(\text{VN} + \text{FP}) \quad (6)$$

El área bajo la curva ROC (AUC) se convierte así en un indicador de la capacidad predictiva del modelo, siendo utilizada para la comparación de diferentes modelos de predicción de susceptibilidad. En este trabajo analizamos separadamente los modelos para las muestras de entrenamiento y validación (Figura 7). Los resultados de los tres métodos son similares, aunque algo mejores para el modelo de RN en la muestra de validación.

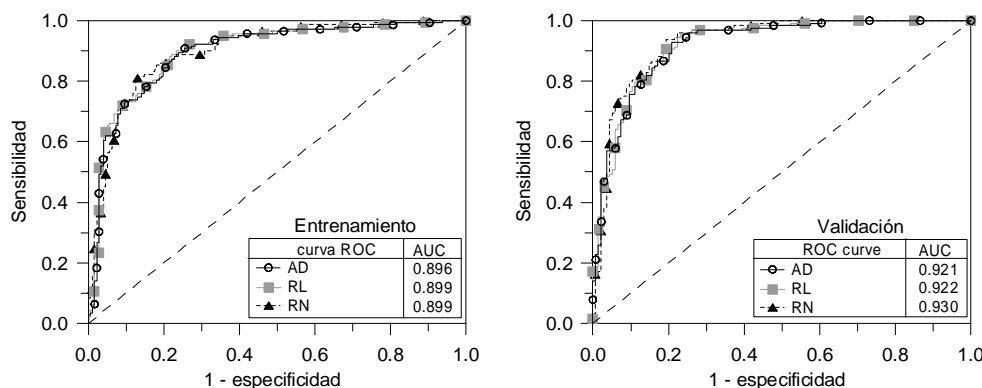


Figura 7. Curvas ROC para muestra de entrenamiento y validación.

Suponiendo que parte de los deslizamientos ocurridos en el pasado representan los deslizamientos futuros, Chung y Fabbri (2003) proponen la validación y verificación de los modelos de predicción de deslizamientos a partir de las curvas *prediction-rate* y *success-rate*. Las curvas de *success-rate* se basan en la comparación entre el modelo de predicción y los deslizamientos utilizados para su creación. Comparando los deslizamientos reservados para validación con el modelo generado se obtienen los estadísticos necesarios para representar la curva *prediction-rate*.

Las curvas (Figura 8) muestran nuevamente unos resultados similares para los tres métodos utilizados en la *success-rate* y ligera ventaja para las RN en la *prediction-rate*.

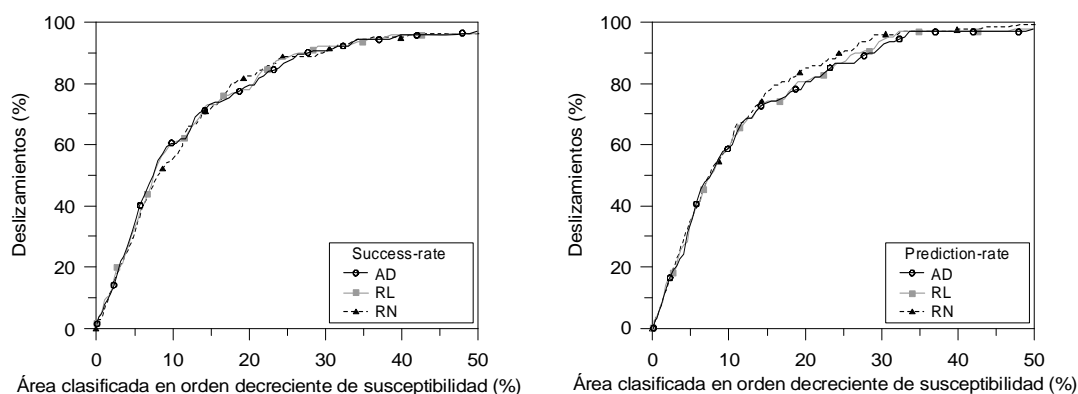


Figura 8: Curvas de éxito y predicción para el 50% de área con mayor susceptibilidad.

6. CONCLUSIONES

- Para una comparación correcta de los métodos de evaluación de la susceptibilidad es necesario que sus resultados sean expresados en términos comparables. En este caso hemos utilizado la probabilidad de que una celda pertenezca al grupo de inestables de acuerdo con la metodología encontrada en análisis discriminante y regresión logística.
- La curva de porcentaje acumulado permiten un análisis complementario de los resultados debido a que representa en curvas separadas las celdas inestables y las celdas totales por nivel de susceptibilidad, al contrario de lo que se observa en el índice de densidad relativa y la curva de tasa de éxito.
- El índice de densidad relativa de deslizamientos muestra que con el uso de cinco intervalos iguales de susceptibilidad, las celdas inestables fueron mejor identificadas con AD. No obstante los histogramas de frecuencias, las curvas ROC, tasa de éxito y tasa de predicción presentaron respuestas ligeramente mejores para la RN.
- En base a la experiencia de los autores, sería más indicado trabajar en una mejora en la calidad de los datos de entrada referentes al inventario de deslizamientos y las variables condicionantes, antes que la búsqueda de metodologías más sofisticadas de evaluación de la susceptibilidad.

REFERENCIAS

- Atkinson, P.M. and Massari, R., 1998. Generalized linear modelling of susceptibility to landsliding in the central Apennines, Italy. *Computer Geoscience*, 24: 373-385.
- Ayalew, L., Yamagishi, H., Marui, H. and Kanno, T., 2005. Landslides in Sado Island of Japan: Part II. GIS-based susceptibility mapping with comparisons of results from two methods and verifications. *Engineering Geology*, 81: 432-445.

- Baeza, C. and Corominas, J., 2001. Assessment of shallow landslide susceptibility by means of multivariate statistical techniques. *Earth Surface Processes. Landforms*, 26: 1251-1263.
- Carrara, A., Crosta, G. and Frattini, P., 2008. Comparing models of debris-flow susceptibility in the alpine environment. *Geomorphology*, 94: 353-378.
- Chung, C.J.F. and Fabbri, A.G., 2003. Validation of Spatial Prediction Models for Landslide Hazard Mapping. *Natural Hazards*, 30: 451-472.
- Demuth, H., Beale, M. and Hagan, M., 2006. *Neural Network Toolbox. For Use with MATLAB. User's Guide Version 5*, The Math Works Inc.
- Dillon, W.R. and Goldstein, M., 1984. *Multivariate analysis. Methods and applications*. John Wiley and Sons. N.Y.
- Duman, T.Y., Can, T., Gokceoglu, C., Nefeslioglu, H.A. and Sonmez, H., 2006. Application of logistic regression for landslide susceptibility zoning of Cekmece Area, Istanbul, Turkey. *Environmental Geology*, 51:241-256.
- Lee, S., Ryu, J.H. and Kim, I.S., 2007. Landslide susceptibility analysis and its verification using likelihood ratio, logistic regression, and artificial neural network models: case study of Youngin, Korea. *Landslides*, 4:327-338.
- Nelles, O., 2000. *Nonlinear system identification: from classical approaches to neural networks and fuzzy models*. Springer, 785pp.
- Ruck, D.W., Rogers, S.K., Kabrisky, M., Oxley, M.E. and Suter, B.W., 1990. The multilayer perceptron as an approximation to a Bayes optimal discriminant function. In: *Neural Networks, IEEE Transactions on*. Volume: 1, Issue: 4, pp. 296-298
- Santacana, N., 2001. *Análisis de susceptibilidad del terreno a la formación de deslizamientos superficiales y grandes deslizamientos mediante el uso de Sistemas de Información Geográfica. Aplicación a la cuenca alta del Río Llobregat*. Tesis Doctoral, Universidad Politécnica de Cataluña.
- Santacana, N., Baeza, C., Corominas, J., de Paz, A. and Marturià, J., 2002. Análisis de la susceptibilidad del terreno a la formación de deslizamientos superficiales mediante el uso de un sistema de información geográfica. Aplicación a la Poble de Lillet (Pirineo Oriental). In: F.J. Ayala-Carcedo y J. Corominas (Editors). *Mapas de susceptibilidad a los movimientos de ladera con técnicas SIG. Fundamentos y aplicaciones en España*. Instituto Geológico y Minero de España. 55-82
- Soeters, R., van Westen, C.J., 1996. Slope instability recognition, analysis and zonation. In Turner&Schuster (Editors), *Landslide Investigation and Mitigation*. Transportation Research Board, Special Report, 247: 129-177, National Academy Press, Washington D.C.
- Visauta, B. and Martori, J., 2003. *Análisis estadístico con SPSS para Windows. Volumen II Estadística multivariante*. McGraw-Hill, Madrid.
- Zêzere, J.L., 2002. Landslide susceptibility assessment considering landslide typology. A case study in the area north of Lisbon (Portugal). *Natural Hazards and Earth System Sciences*, 2: 73-82.